

# **Module 10**

## **Principal Component Analysis / Shrinkage**

Kentaro Nakamura

GOV 2003

April 17th, 2026

# Agenda

- Review of Midterm
  
- Principal Component Analysis (PCA)
  
- Shrinkage method
  - Ridge regression
  - Lasso regression

# Midterm: Sensitivity Analysis (1)

- Consider the following data generating process:

$$\mathbf{Y} = \beta \mathbf{T} + \mathbf{X}\boldsymbol{\gamma} + \delta \mathbf{U} + \boldsymbol{\epsilon} \quad \mathbb{E}[\boldsymbol{\epsilon}_i \mid \mathbf{T}, \mathbf{X}, \mathbf{U}] = 0$$

- $\mathbf{U}$  is unobserved
- Now, suppose that you fit the regression above without  $\mathbf{U}$ . How does the bias look like?
- Now, using FWL theorem,

$$\hat{\beta} = (\mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{T})^{-1} \mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{Y} = \frac{\mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{Y}}{\mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{T}}$$

## Midterm: Sensitivity Analysis (2)

- As we have  $\mathbf{Y} = \beta \mathbf{T} + \mathbf{X}\gamma + \delta \mathbf{U} + \epsilon$ ,

$$\begin{aligned}\hat{\beta} &= \frac{\mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{Y}}{\mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{T}} \\ &= \beta + \frac{\mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{X}}{\mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{T}} \gamma + \delta \frac{\mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{U}}{\mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{T}} + \frac{\mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \epsilon}{\mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{T}}.\end{aligned}$$

- By the weak law of large numbers,

$$\begin{aligned}\frac{1}{n} \mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \epsilon &\xrightarrow{p} \mathbb{E}[T_1 (I - \mathbf{P}_{X_1}) \epsilon_1] \\ &= \mathbb{E}[T_1 (I - \mathbf{P}_{X_1})] \mathbb{E}[\epsilon_i | \mathbf{T}, \mathbf{X}, \mathbf{U}] = 0,\end{aligned}$$

- Similarly, the second term is

$$\frac{1}{n} \mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{X} = 0.$$

and thus it is also negligible.

## Midterm: Sensitivity Analysis (3)

- Finally, the middle term is written as

$$\delta \frac{\mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{U}}{\mathbf{T}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{T}} = \delta \times \frac{\frac{1}{n} \{(\mathbf{I}_n - \mathbf{P}_X) \mathbf{T}\}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{U}}{\frac{1}{n} \{(\mathbf{I}_n - \mathbf{P}_X) \mathbf{T}\}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{T}}.$$

- The numerator is

$$\frac{1}{n} \{(\mathbf{I}_n - \mathbf{P}_X) \mathbf{T}\}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{U} \xrightarrow{p} \mathbb{E}[T_i^{\perp \mathbf{X}} U_i^{\perp \mathbf{X}}] = \text{Cov}(T_i^{\perp \mathbf{X}}, U_i^{\perp \mathbf{X}})$$

- The denominator is

$$\frac{1}{n} \{(\mathbf{I}_n - \mathbf{P}_X) \mathbf{T}\}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{T} \xrightarrow{p} \mathbb{E}[(T_i^{\perp \mathbf{X}})^2] = \mathbb{V}[T_i^{\perp \mathbf{X}}]$$

- Therefore, we get

$$\hat{\beta} \xrightarrow{p} \beta + \delta \times \frac{\text{Cov}(T_i^{\perp \mathbf{X}}, U_i^{\perp \mathbf{X}})}{\mathbb{V}[T_i^{\perp \mathbf{X}}]}$$

## Midterm: Sensitivity Analysis (4)

- Now,  $\delta$  is the coefficient of unobserved confounder. So, by FWL,

$$\delta = \frac{\text{Cov}(Y_i^{\perp \mathbf{X}, T}, U_i^{\perp \mathbf{X}, T})}{\mathbb{V}[U_i^{\perp \mathbf{X}, T}]}$$

- Also, recall that

$$\mathbb{V}(Y_i - \hat{Y}_i) = (1 - R^2)\mathbb{V}(Y_i)$$

- This is because  $\mathbb{V}(Y_i) = \mathbb{V}(\hat{Y}_i) + \mathbb{V}(Y_i - \hat{Y}_i)$  and  $R^2 = \frac{\mathbb{V}(\hat{Y}_i)}{\mathbb{V}(Y_i)}$

- Thus,

$$\mathbb{V}[U_i^{\perp \mathbf{X}, T}] = \mathbb{V}[U_i^{\perp \mathbf{X}}] \times (1 - R_{U \sim T | \mathbf{X}}^2) = \mathbb{V}[U_i^{\perp \mathbf{X}}] \times (1 - R_{T \sim U | \mathbf{X}}^2)$$

and we get

$$|\hat{\delta}| = \sqrt{R_{Y \sim U | T, \mathbf{X}}^2} \sqrt{\frac{\hat{\mathbb{V}}(Y^{\perp \mathbf{X}, T})}{\hat{\mathbb{V}}(U^{\perp \mathbf{X}})(1 - R_{T \sim U | \mathbf{X}}^2)}}.$$

## Midterm: Sensitivity Analysis (5)

- Finally, notice that

$$R_{T \sim U | \mathbf{X}}^2 = \text{Corr}(T_i^\perp \mathbf{X}, U_i^\perp \mathbf{X})^2 = \frac{\text{Cov}(T_i^\perp \mathbf{X}, U_i^\perp \mathbf{X})^2}{\mathbb{V}[T_i^\perp \mathbf{X}] \mathbb{V}[U_i^\perp \mathbf{X}]}$$

and thus

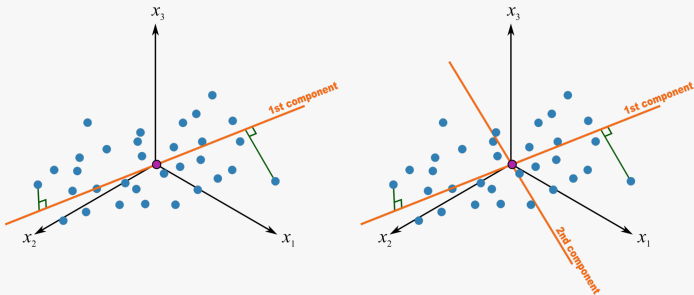
$$\left| \frac{\text{Cov}(T_i^\perp \mathbf{X}, U_i^\perp \mathbf{X})}{\mathbb{V}[T_i^\perp \mathbf{X}]} \right| = \sqrt{R_{T \sim U | \mathbf{X}}^2} \times \sqrt{\frac{\mathbb{V}[U_i^\perp \mathbf{X}]}{\mathbb{V}[T_i^\perp \mathbf{X}]}}$$

- Therefore,

$$\begin{aligned} |\text{bias}| &= |\hat{\delta}| \times \left| \frac{\text{Cov}(T_i^\perp \mathbf{X}, U_i^\perp \mathbf{X})}{\mathbb{V}[T_i^\perp \mathbf{X}]} \right| \\ &= \sqrt{R_{Y \sim U | T, \mathbf{X}}^2} \sqrt{\frac{\mathbb{V}[Y_i^\perp \mathbf{X}, T]}{\mathbb{V}[U_i^\perp \mathbf{X}] \times (1 - R_{T \sim U | \mathbf{X}}^2)}} \sqrt{R_{T \sim U | \mathbf{X}}^2} \times \sqrt{\frac{\mathbb{V}[U_i^\perp \mathbf{X}]}{\mathbb{V}[T_i^\perp \mathbf{X}]}} \\ &= \sqrt{R_{Y \sim U | T, \mathbf{X}}^2} \times \sqrt{\frac{R_{T \sim U | \mathbf{X}}^2}{1 - R_{T \sim U | \mathbf{X}}^2}} \times \sqrt{\frac{\mathbb{V}[Y_i^\perp \mathbf{X}, T]}{\mathbb{V}[T_i^\perp \mathbf{X}]}} \end{aligned}$$

# Principal Component Analysis (PCA)

- **Goal:** Dimension reduction
  - We often have so many variables, and we want to reduce the dimensionality and discover which dimension matters the most
  - We want to transform high-dimensional, correlated data, into a smaller set of uncorrelated variable (known as **principal component**)



- Notice that the first principal component is the direction where the variance is maximized.

# Principal Component Analysis (PCA) (1)

- Let  $x_1, \dots, x_N \in \mathbb{R}^p$  be the observations, and let  $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$ 
  - We can then define the centered data matrix

$$\bar{X} = \begin{bmatrix} (x_1 - \bar{x})^\top \\ \vdots \\ (x_N - \bar{x})^\top \end{bmatrix} \in \mathbb{R}^{N \times p}$$

- The goal is that to find  $j$ -th principal component ( $j \in \{1, \dots, J\}$ ) that is the linear transformation of the data

$$s_j = \bar{X} a_j \in \mathbb{R}^N$$

where  $a_j$  is the coefficient

- We want this transformation to satisfy (i) length is 1 and (ii) variance is maximized

## Principal Component Analysis (PCA) (2)

- Let's calculate the variance of  $s_j$ . Now, as  $s_j$ 's average is 0,

$$\text{Var}[s_j] = \frac{1}{N} s_j^\top s_j = \frac{1}{N} a_j^\top \bar{X}^\top \bar{X} a_j = a_j^\top \text{Var}[\bar{X}] a_j$$

- We want to find the projection vector  $a_j$  such that (i) length is 1 and (ii) variance is maximized. That optimization is written as

$$\max a_j^\top \widehat{\text{Var}[\bar{X}]} a_j \quad \text{such that} \quad a_j^\top a_j = 1$$

## Principal Component Analysis (PCA) (3)

- By using Lagrangian multiplier, we have the Lagrangian

$$\mathcal{L}(a_j, \lambda) = a_j^\top \widehat{\text{Var}[\bar{X}]} a_j - \lambda_j (a_j^\top a_j - 1)$$

and the derivative is

$$\partial_{a_j} \mathcal{L}(a_j, \lambda_j) = 2 \widehat{\text{Var}[\bar{X}]} a_j - 2 \lambda_j a_j = 0$$

$$\partial_\lambda \mathcal{L}(a_j, \lambda_j) = a_j^\top a_j - 1 = 0$$

- Now notice that the first constraint is

$$\widehat{\text{Var}[\bar{X}]} a_j = \lambda_j a_j$$

which is eigenvalue problem!

- Lagrangian multiplier  $\lambda_j$  becomes eigenvalue
- Note that we need to scale the eigenvector so that its length becomes 1.

## Principal Component Analysis (PCA) (4)

- As a result, PCA corresponds to solving the eigenvalue decomposition of variance matrix  $\widehat{\text{Var}}[\bar{X}]$ ;
  - $\widehat{\text{Var}}[\bar{X}] = \frac{1}{N} \bar{X}^\top \bar{X}$  is symmetric variance, so spectral theorem guarantees the existence of EVD
- Moreover, the variance of  $j$ -th principal component  $s_j$  is

$$\widehat{\text{Var}}[s_j] = \frac{1}{N} (\bar{X} a_j)^\top (\bar{X} a_j) = a_j^\top \widehat{\text{Var}}[\bar{X}] a_j$$

and thus by eigenvalue decomposition  $\widehat{\text{Var}}[\bar{X}] = A \Lambda A^\top$  with  $A = [a_1, \dots, a_p]$  and  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$ , so

$$\widehat{\text{Var}}[s_j] = a_j^\top A \Lambda A^\top a_j = e_j^\top \Lambda e_j = \lambda_j$$

- That is, the eigenvalue corresponds to the variance of principal component

# Principal Component Analysis (PCA) (5)

- Define the total variance of centered regressor as

$$\text{Total variance} := \sum_{j=1}^p \text{Var}(X_j)$$

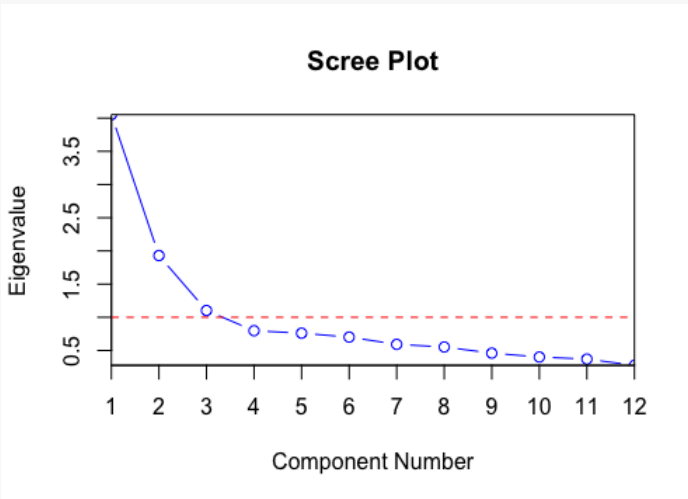
- Now, notice that

$$\text{Total variance} = \text{tr}(\widehat{\text{Var}}[\bar{X}]) = \text{tr}(A\Lambda A^\top) = \text{tr}(\Lambda) = \sum_{j=1}^p \lambda_j$$

- How many principal components should we use?
  - If we use all  $p$  components (same number as the number of variables), we do not reduce the dimensions, and we can explain all the variations in the data
  - However, if we only use the very small number of principal components, we might not explain the variation in the data well.

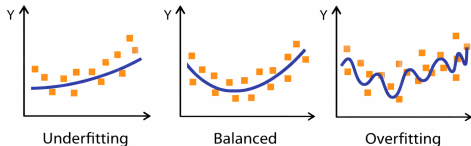
## Principal Component Analysis (PCA) (6)

- Since eigenvalue equals the variance of the principal component, we can decide how many components we would use from the data
  - You can see that at which point adding the dimensions do not help explaining the variance in the data.



# Regularization: Motivation (1)

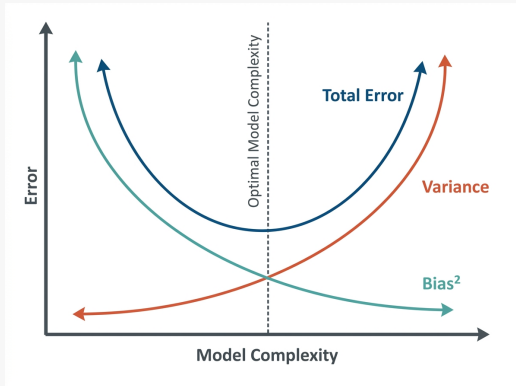
- We want to make the function flexible enough to minimize the bias
- However, we do not want the overfitting



- As figure illustrates, overfitting is caused by high variance
  - There is a trade-off: underfitting happens because of restrictive model, whereas overfitting happens because of complex model

## Regularization: Motivation (2)

- Last time, we learn that to minimize the out-of-sample prediction error, we want to minimize both bias and variance
  - Only minimizing bias might not be optimal for variance



- Regularization introduces bias by limiting the model complexities, but shrinks variance

## Ridge regression (1)

- Think about constraining the model complexity for linear model by adding constraints on coefficients

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{\alpha, \beta} \sum_{i=1}^n (Y_i - \alpha - \mathbf{x}_i^T \beta)^2 \quad \text{s.t.} \quad \|\beta\|_2^2 \leq t$$

- Note that  $\|\beta\|_2^2 = \sum_{j=1}^p \beta_j^2$
- Using Lagrangian multiplier, the dual problem is written as

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{\alpha, \beta} \sum_{i=1}^n (Y_i - \alpha - \mathbf{x}_i^T \beta)^2 + \lambda (\|\beta\|_2^2 - t)$$

with  $\lambda \geq 0$

- For a fixed  $\lambda$ , the term  $-\lambda t$  is constant, so it is equal to minimizing<sup>1</sup>

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{\alpha, \beta} \sum_{i=1}^n (Y_i - \alpha - \mathbf{x}_i^T \beta)^2 + \lambda \|\beta\|_2^2$$

---

<sup>1</sup>Notice that we do not optimize the Lagrangian multiplier  $\lambda$ . So, the equivalence here is that there exists a correspondence between  $\lambda$  and  $t$  and this is one-to-one, but the mapping is specific to each data set.

## Ridge regression (2)

- Now, let  $\tilde{Y}_i = Y_i - \bar{Y}$  and  $\tilde{\mathbf{X}}_i = \mathbf{X}_i - \bar{\mathbf{X}}$ , and the model becomes

$$\tilde{Y}_i = \tilde{\mathbf{X}}_i^\top \boldsymbol{\beta} + \epsilon_i$$

- By standardization, we do not have an intercept
- Then, optimization problem is

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \|\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2$$

- By taking the derivative, we get

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\beta}} \left( \|\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2 \right) \\ = -2\tilde{\mathbf{X}}^\top (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}) + 2\lambda\boldsymbol{\beta} = 0 \end{aligned}$$

- This gives us the closed-form solution

$$\hat{\boldsymbol{\beta}} = (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} + \lambda \mathbf{I})^{-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{Y}}$$

## Ridge regression (3)

- When  $\lambda > 0$ ,  $\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} + \lambda \mathbf{I}$  is always invertible
- Recall that matrix is invertible if it is positive definite
  - Now, consider the quadratic form

$$\mathbf{v}^\top (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} + \lambda \mathbf{I}) \mathbf{v} = \|\tilde{\mathbf{X}}^\top \mathbf{v}\|_2^2 + \lambda \|\mathbf{v}\|_2^2$$

- Now, even if  $\|\tilde{\mathbf{X}}^\top \mathbf{v}\|_2^2 = 0$ ,  $\|\mathbf{v}\|_2^2 > 0$  as  $\mathbf{v} \neq \mathbf{0}$
- So, when  $\lambda > 0$ ,  $\mathbf{v}^\top (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} + \lambda \mathbf{I}) \mathbf{v} > 0$ , meaning that  $\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} + \lambda \mathbf{I}$  is always invertible
- This is different from OLS
  - In OLS,  $\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}$  can be not invertible when regressors are rank-deficient (or known as perfect collinearity)
  - This is likely when you have so many regressors
  - Ridge regression solves this problem

## Ridge regression (4)

- However, it is important to recognize that the coefficient of ridge regression is biased!
- To see why,

$$\begin{aligned}\mathbb{E}[\hat{\beta} \mid \mathbf{X}] &= (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} + \lambda \mathbf{I})^{-1} \tilde{\mathbf{X}}^\top \mathbb{E} \tilde{\mathbf{Y}} \mid \mathbf{X}] \\ &= (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} + \lambda \mathbf{I})^{-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \beta\end{aligned}$$

which is different from  $\beta$  unless  $\lambda = 0$

- This makes sense because regularization makes the variance smaller but introduces the bias (called regularization bias)
  - If you need to use models with regularization for statistical inference, you need to make your estimation robust to regularization bias
  - This is exactly what double machine learning is doing.

## Ridge regression (5)

- Recall that SVD gives you  $\tilde{\mathbf{X}} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ 
  - Here,  $\mathbf{U} \in \mathbb{R}^{n \times r}$ ,  $\mathbf{V} \in \mathbb{R}^{p \times r}$  and  $\mathbf{D} = \text{diag}(d_1, \dots, d_r)$  where  $r = \text{rank}(\tilde{\mathbf{X}})$ .
- Plugging in the closed-form solution for ridge, we get

$$\begin{aligned}\hat{\beta} &= (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} + \lambda \mathbf{I})^{-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{Y}} \\ &= (\mathbf{V}\mathbf{D}^2\mathbf{V}^\top + \lambda \mathbf{I})^{-1} \mathbf{V}\mathbf{D}\mathbf{U}^\top \tilde{\mathbf{Y}} \\ &= \mathbf{V}(\mathbf{D}^2 + \lambda \mathbf{I})\mathbf{V}^\top \mathbf{V}\mathbf{D}\mathbf{U}^\top \tilde{\mathbf{Y}}\end{aligned}$$

where the last equality is because  $\mathbf{V}$  is orthonormal.

- Thus, let  $\mathbf{z} = \mathbf{U}^\top \tilde{\mathbf{Y}}$ . Then,

$$\hat{\beta} = \sum_{j=1}^r \frac{d_j}{d_j^2 + \lambda} z_j \mathbf{v}_j$$

- Notice that this formulation does not require you to calculate the inverse (i.e., computationally efficient)
- Similar logic to QR decomposition to OLS

## Added: Bias-Variance Tradeoff of Ridge

- Recall that we derived

$$\hat{\beta} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I})^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{Y}}$$

- For the sake of simplicity, homoskedasticity  $\mathbb{V}[\epsilon | \mathbf{X}] = \sigma^2 \mathbf{I}_n$
- Then, the variance is written as

$$\begin{aligned}\mathbb{V}(\hat{\beta} | \tilde{\mathbf{X}}) &= \mathbb{V}((\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I})^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{Y}} | \tilde{\mathbf{X}}) \\ &= \mathbb{V}((\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I})^{-1} \tilde{\mathbf{X}}^T [\tilde{\mathbf{X}}\beta + \epsilon] | \tilde{\mathbf{X}}) \\ &= \mathbb{V}((\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I})^{-1} \tilde{\mathbf{X}}^T \epsilon | \tilde{\mathbf{X}}) \\ &= (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I})^{-1} \tilde{\mathbf{X}}^T \mathbb{V}(\epsilon | \tilde{\mathbf{X}}) \tilde{\mathbf{X}} (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I})^{-T} \\ &= \sigma^2 (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I})^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I})^{-T}\end{aligned}$$

- Notice that as  $\lambda$  increases, the variance decreases
  - This is the trade-off between bias and variance in ridge regression

# Lasso regression (1)

- Instead of using L2 norm  $\|\beta\|_2^2$  for penalty, consider

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{\alpha, \beta} \sum_{i=1}^n (Y_i - \alpha - \mathbf{X}_i^\top \beta)^2 + \lambda \|\beta\|_1$$

then  $\|\beta\|_1 = \sum_{j=1}^p |\beta_j|$ .

- As with the previous case, the formulation above is dual problem of

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{\alpha, \beta} \sum_{i=1}^n (Y_i - \alpha - \mathbf{X}_i^\top \beta)^2 \quad \text{s.t.} \quad \|\beta\|_1 \leq t$$

## Lasso regression (2)

- Lasso regression makes the coefficient tend to be 0, whereas Ridge regression makes the coefficient smaller
  - Intuition can be seen from this figure (try to understand!)

